# GHAJAR
# EXHIBIT 21

1          UNITED STATES DISTRICT COURT

2          NORTHERN DISTRICT OF CALIFORNIA

3          SAN FRANCISCO DIVISION

4

5   RICHARD KADREY, AN              : CASE NO.

6   INDIVIDUAL; SARAH SILVERMAN,    : 3:23-cv-03417-VC

7   AN INDIVIDUAL; CHRISTOPHER      :

8   GOLDEN, AN INDIVIDUAL,          :

9          PLAINTIFFS               :

10  VS.                             :

11  META PLATFORMS, INC., A         :

12  DELAWARE CORPORATION,           :

13          DEFENDANT               :
    _____:

14       HIGHLY CONFIDENTIAL - ATTORNEYS' EYES ONLY

15

16     VIDEOTAPED DEPOSITION OF JONATHAN KREIN, PH.D.

17          SAN FRANCISCO, CALIFORNIA

18          THURSDAY, MARCH 6, 2025

19

20   REPORTED BY:

21   DEBBIE LEONARD, CSR, RDR, CRR

22   CSR NO. 14350

23   JOB No. 7189213

24

25   PAGES 1 to 113

                                        Page 1

HIGHLY CONFIDENTIAL - ATTORNEYS' EYES ONLY

1          MR. YOUNG:  Objection.  Form.

2          THE WITNESS:  So, for example, in paragraph 81,

3     I discuss deduplication, lay it out there.

4          So the substance of the opinions that I've

5     formed with respect to deduplication is around its impact

6     on the generation of copies.

7          Exactly what this particular instance of

8     deduplication was used for by Meta was not something I

9     was asked to consider.

10    BY MR. WEINSTEIN:

11       Q    Understood.

12            If you could turn to paragraph 97.  You

13    identified as part of your analysis certain Python

14    scripts that would remove certain information from the

15    content in The Pile dataset, correct?

16       A    Could you repeat the question, please?

17       Q    Yes, sir.

18            You identified in paragraph 97 a Python script

19    that removes information from various aspects of the

20    The Pile dataset, correct?

21       A    Yes.

22       Q    Okay.  Now, the Python script that you

23    identified, pile_clean_v0.py, that is not limited to

24    Books3; it also applies to ArXiv, A-R-X-I-V; GitHub; and

25    other portions of Books3, correct?

Page 77

1    A    Correct.

2    Q    Okay.  And then let's go to paragraph 99.  And

3  is this a list of key words that are removed by one of

4  the Python scripts that you analyzed with respect to

5  The Pile?

6    A    So this list of terms refers to paragraph 98 of

7  my report, which discusses a script called

8  pile_clean_v1.py, and it is referring to terms that were

9  used to identify lines or segments of -- of The Pile data

10  to be removed.

11    Q    And you have a list on paragraph 99 of key words

12  that are used for identifying [indiscernible], correct?

13         (Court reporter clarification.)

14  BY MR. WEINSTEIN:

15    Q    Paragraph 99, you have a list of key words for

16  use in identifying material to be removed, correct?

17         MR. YOUNG:  Objection.  Form.

18         THE WITNESS:  So as I just mentioned, the list

19  in paragraph 99 refers to the script in paragraph 98, and

20  these are key words that are for use in identifying

21  material to be removed, yes.

22  BY MR. WEINSTEIN:

23    Q    Okay.  And they include "Published by,"

24  "Produced by," "Copyright," "Cover copyright," "Cover

25  designed by," "Edited by," "Prologue," "PREFACE,"

Page 78

1    "Epilogue," "Acknowledgments," "About the Author,"

2    "Slashwords [sic] Edition," "All rights reserved," "Title

3    Page," "Table of Contents," "Introduction," "Cover,"

4    "Title Page," and again the word "Copyright."

5         Do you see that?

6    A    I do see that, with the exception that I believe

7    you may have said "Slashwords Edition," and it says

8    "Smashwords Edition."

9    Q    Thank you.

10        MR. WEINSTEIN:   What's next in order?

11        THE REPORTER:   Number 5.

12        MR. WEINSTEIN:   5.   Thank you.

13        I'd like to mark as Exhibit 5 a copy of a

14   document entitled pile_clean_v1.py.   That's

15   META-KADREY-SC-000625 through 629.

16        (Krein Exhibit 5 marked for identification.)

17   BY MR. WEINSTEIN:

18   Q    Dr. Krein, the court reporter just handed you a

19   document that's called pile_clean_v1.py.   This is one of

20   the source code files that you looked at in connection

21   with your opinions related to The Pile, correct?

22   A    Correct.

23   Q    And this is a Python file, correct?

24   A    Yes.

25   Q    And just humor me for the record.   What does

1    Python refer to?

2        A    Python is a programming language.

3        Q    If you could turn to the page that begins 000628

4    and look around line 304.

5        A    I'm on page 628.  Well --

6        Q    And you see line 304?

7        A    I do, yes.

8        Q    And that's the main program, correct?

9        A    That's the main entry point.

10        Q    Correct.  So when the script is executed, that's

11    the first part of the code that will get executed,

12    correct?

13        A    No.

14        Q    Okay.  Will it get executed at all?

15        A    Yeah.  To clarify, on the next page, you'll see

16    at line 371, that's where main is called from.  And then

17    at 304, it executes.

18        Q    I see.  So main will then call main, and then

19    the trip will function, right?

20        A    Yeah.

21        Q    Okay.  And do you see where it says "cleaners

22    equals," and there's a list?

23        A    I do, yes.

24        Q    What does that list refer to?

25        A    So my understanding is that these would be parts

1    of -- parts of The Pile dataset.

2        Q    Are these portions of The Pile dataset for which

3    cleaning operations are performed?

4            MR. YOUNG:  Objection.  Form.

5            THE WITNESS:  So these are portions, as Meta has

6    named them here, of The Pile dataset for which there are,

7    yes, cleaning operations that get performed on any part

8    of the data that's designated under one of these

9    designations here.

10   BY MR. WEINSTEIN:

11       Q    And there are 22 of them, correct,

12   approximately?

13       A    There may be.  I'd have to go through and count.

14       Q    Okay.  And you see that on line 313, it says

15   "Books3"?

16           Do you see that?

17       A    Line 313 refers to "Books3."

18       Q    Okay.  Is that the cleaner that relates to

19   Books3?

20           MR. YOUNG:  Objection.  Form.

21   BY MR. WEINSTEIN:

22       Q    Or it's identifying the cleaner that relates to

23   Books3?

24       A    So that is a text designation that points to the

25   name of a function which is called "book3" at line 63.

Page 81

1        When the parameter "Books3" is provided, then it

2    results in the function at 63 being executed.

3        Q     Understood.

4              We can go back to your report now, sir,

5    Exhibit 1.   And if you could turn to paragraph 100, this

6    is another file you identified, a .scala file, that

7    removes data from the Books3 dataset, correct?

8        A     So paragraph 100 of my report:

9    "BooksHeaderCleaner.scala removes copyright information

10   from the Books3 dataset."

11       Q     And just so we know, what is a Scala file?

12       A     It is a type of programming file.

13       Q     And then you have a long list of key words that

14   goes from page 46 to 47 in your report.   Are these key

15   words that are removed from the Books3 dataset when this

16   program is executed?

17       A     So the key words are used to identify segments

18   or lines -- I believe in this case, it's -- it may be

19   lines; I'd have to look back -- but segments of text to

20   be removed from the data.

21       Q     And at the top of page 47, you mention that

22   there are a number of other terms, like the word

23   "facebook" being removed?

24       A     The word "facebook" is one of the terms that is

25   additionally used in filtering lines.

1    Q    And "www.," correct?

2    A    That also is one of the terms that I've listed

3    here in paragraph 100.

4    Q    And "notebook" is another one, correct?

5    A    "Notebook" as a single word, and "note-book" as

6    well.

7    Q    Okay.  Do you know why this script would be

8    designed to remove the word "facebook" from the original

9    text?

10         MR. YOUNG:  Objection to form.  Scope.

11         THE WITNESS:  I did not examine that particular

12   question and didn't form any opinion about it.

13   BY MR. WEINSTEIN:

14   Q    Okay.  And then on paragraph 101, you discuss

15   another source code file called BooksFooterCleaner.scala,

16   and you have another list of key words, correct?

17   A    Yes, paragraph 101 does refer to another

18   programming script that also lists a set of key words

19   that are used to identify text.

20   Q    And these include things like "already a

21   subscriber," "thank you for downloading," "thank you for

22   purchasing," among others, correct?

23   A    To the extent that you read those out of this

24   list, then they would, of course, be in the list.

25   Q    Now, when you analyzed -- withdrawn.

Page 83

1                   C E R T I F I C A T E

2

3          I, Debbie Leonard, Certified Shorthand Reporter

4     No. 14350 for the State of California, do hereby

5     certify:

6          That the foregoing deposition was taken before me

7     at the time and place therein set forth, at which time

8     the witness was put under oath by me; that the testimony

9     of the witness and all objections made at the time of the

10    examination were recorded stenographically by me, were

11    thereafter transcribed by me by means of computer; and

12    that the foregoing is a true record of same.

13         I further certify that I am neither counsel for

14    nor related to any party to said action, nor in any way

15    interested in the outcome thereof.

16         IN WITNESS WHEREOF, I have subscribed my name

17    this 13th day of March, 2025.

18

19

20

21          _Debbie Leonard_

22          Debbie Leonard, CSR, RDR, CRR

23          CSR NO. 14350

24

25